

Introducción a la pan-genómica microbiana

IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana,
TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>

Microbial genome evolution - the pan-genome

Pablo Vinuesa
Centro de Ciencias Genómicas – UNAM
vinuesa @ ccg_unam_mx
<http://www.ccg.unam.mx/~vinuesa/>
 @pvinmex

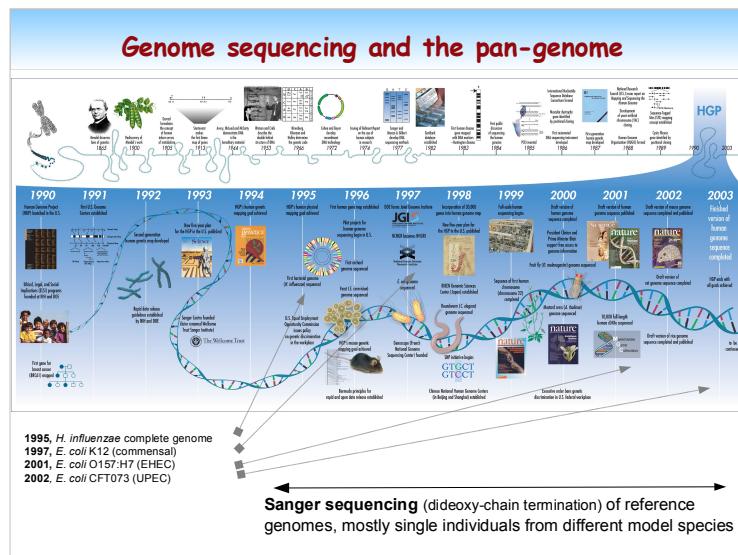
Introducción a la Filoinformática:
Pan-genómica y Filogenómica microbiana –
NNB & CCG-UNAM, 1-5 Agosto 2022
<https://github.com/vinuesa/TIB-filoinfo>



Genome evolution and the pan-genome

Contents

1. Basic concepts in bacterial comparative genomics – defining the pan-genome
 - genome mosaicism and the pan-genome
 - Structure of the pan-genome: the core and accessory/flexible (shell and cloud) genomes
 - homologous recombination and the genetic structure of bacteria populations
 - MGEs (ISs, GIs, Integrons, ICEs, prophages and plasmids) and the accessory genome
2. Computational challenges and approaches in pan-genome research
 - computational approaches to identify homologs: BDBH, COGs, OMCL
3. Popular software implementations to compute and analyze bacterial pan-genomes
 - 3.1 Computationally-intensive methods for accurate homology inference in diverged species
 - get_homologues
 - 3.2 Heuristic clustering methods for rapid species-level pan-genome analysis
 - Roary
4. Phylogenomic approaches for pan-genome analyses
 - 4.1 Computationally-intensive approaches:
 - get_phylomarkers maximum-likelihood core-genome phylogenies with optimal markers
 - get_phylomarkers and ML or parsimony pan-genome phylogenies
 - 4.2 Heuristic approaches to resolve population-level phylogenies
 - clustering pan-genome presence-absence distance matrix



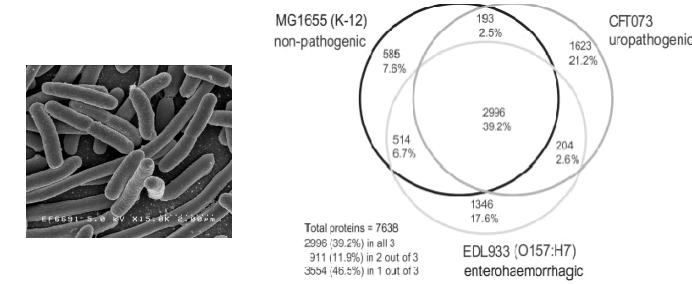
Basic concepts in bacterial comparative genomics – Mosaic genome structure and the pan-genome

Extensive mosaic structure revealed by the complete genome sequence of uropathogenic *Escherichia coli*

R. A. Welch*, V. Burland†‡, G. Plunkett III*, P. Redford*, P. Roesch*, D. Rasko§, E. L. Buckles§, S.-R. Liou†, A. Boutin**
J. Hackett†‡, D. Strout‡, G. F. Mayhew‡, D. J. Rose‡, S. Zhou†‡, D. C. Schwartz†‡, N. T. Perna§§, H. L. T. Mobley‡,
M. S. Donnenberg§, and F. R. Blattner†

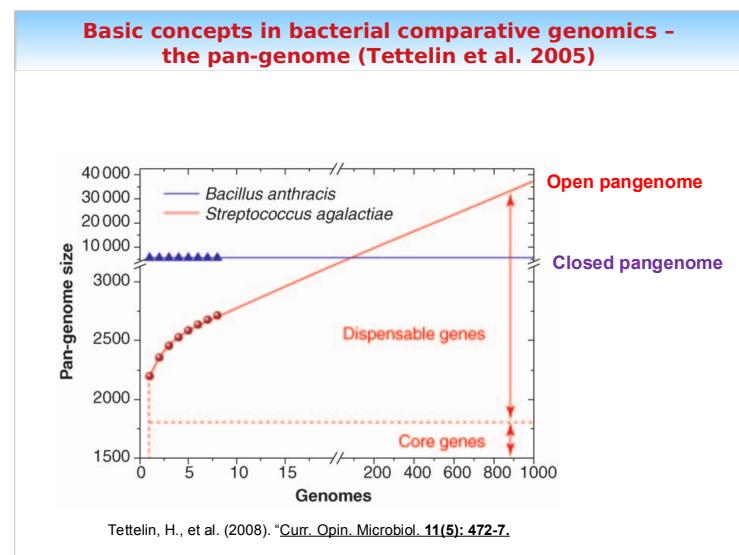
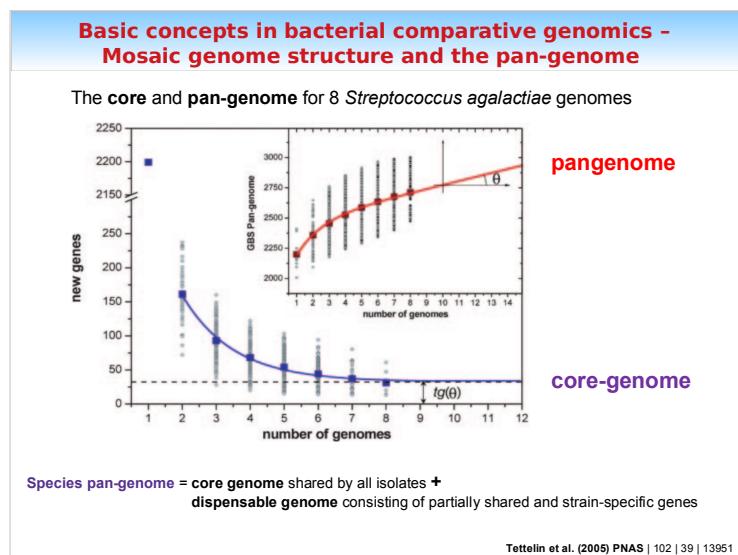
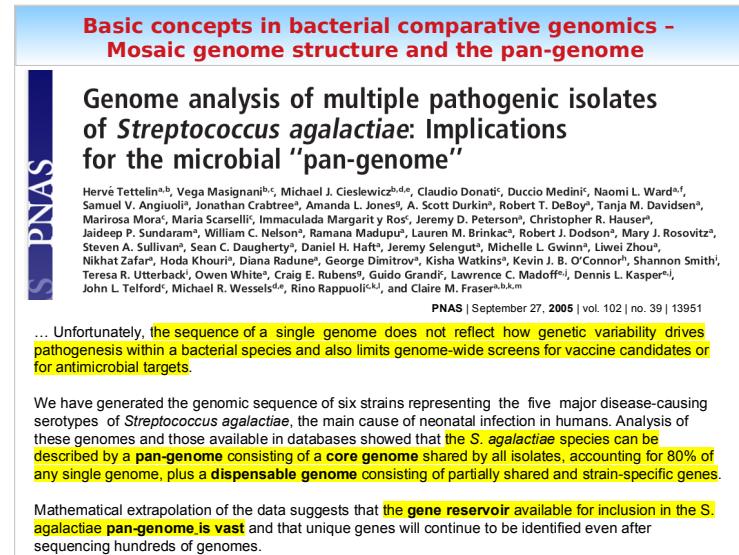
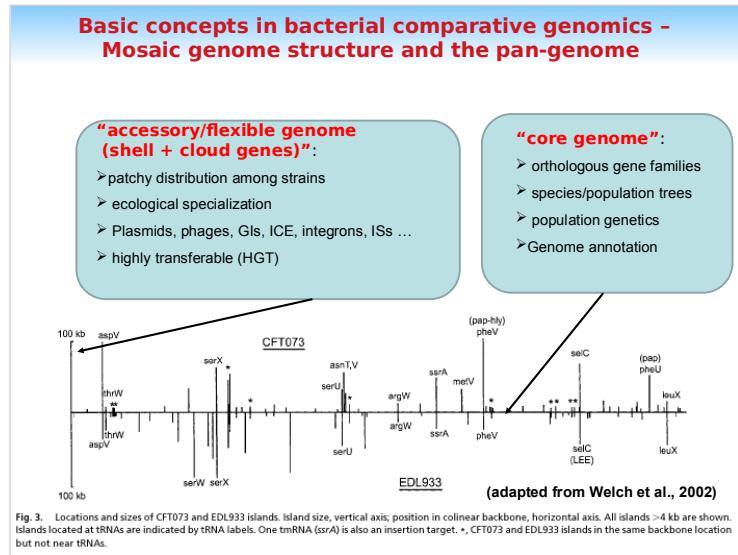
17020–17024 | PNAS | December 24, 2002 | vol. 99 | no. 26

Edited by John J. Mekalanos, Harvard Medical School, Boston, MA, and approved October 24, 2002 (received for review August 30, 2002)



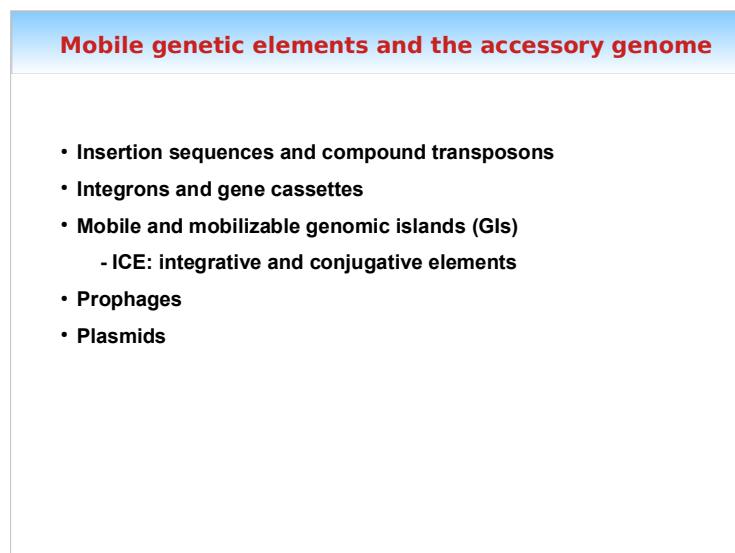
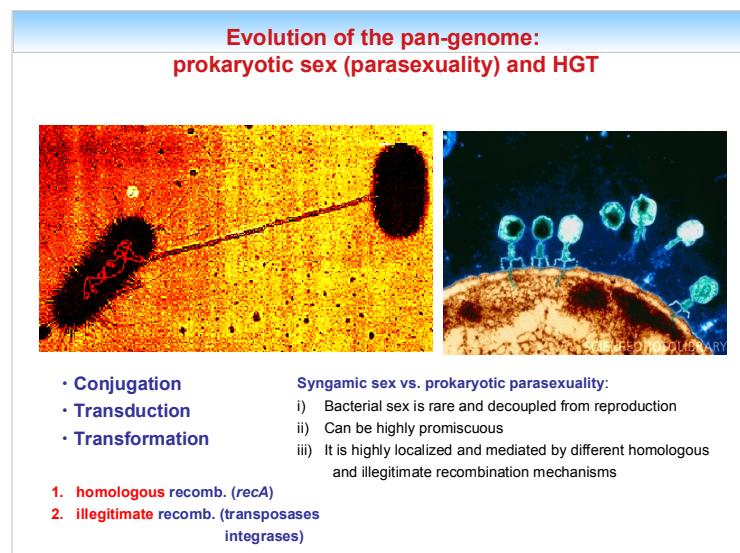
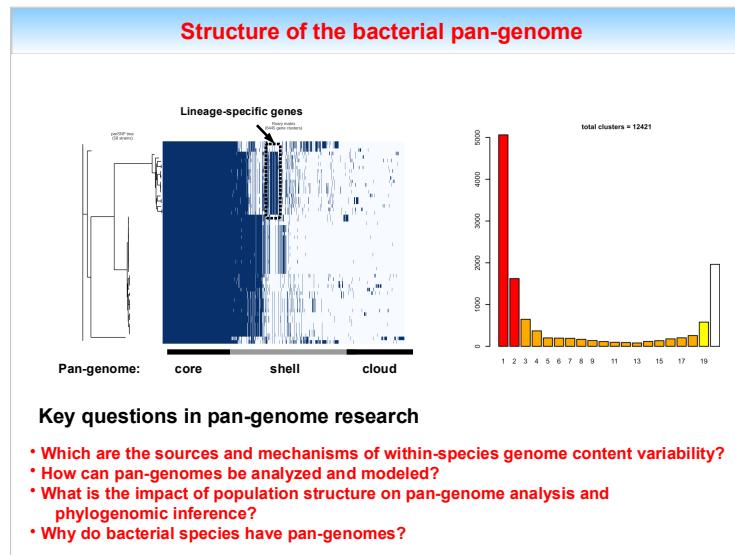
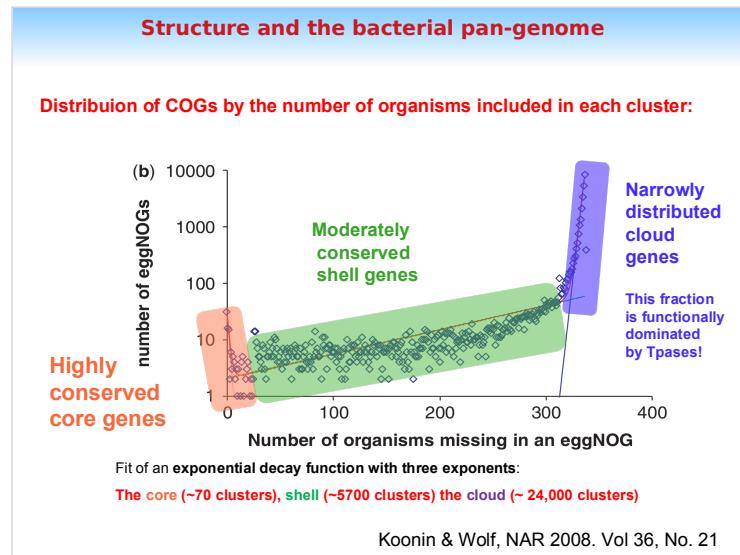
Introducción a la pan-genómica microbiana

IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana, TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



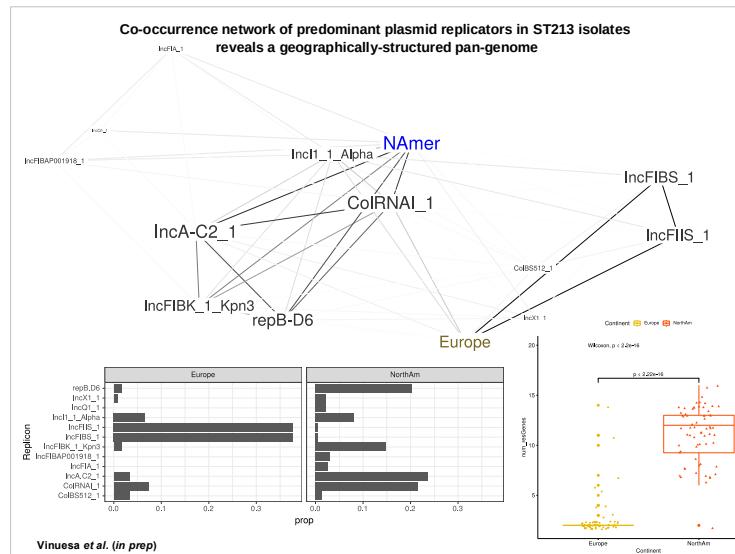
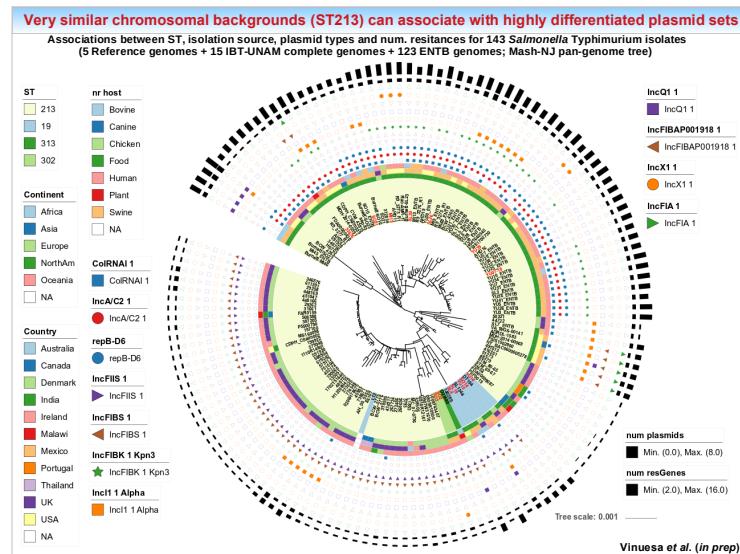
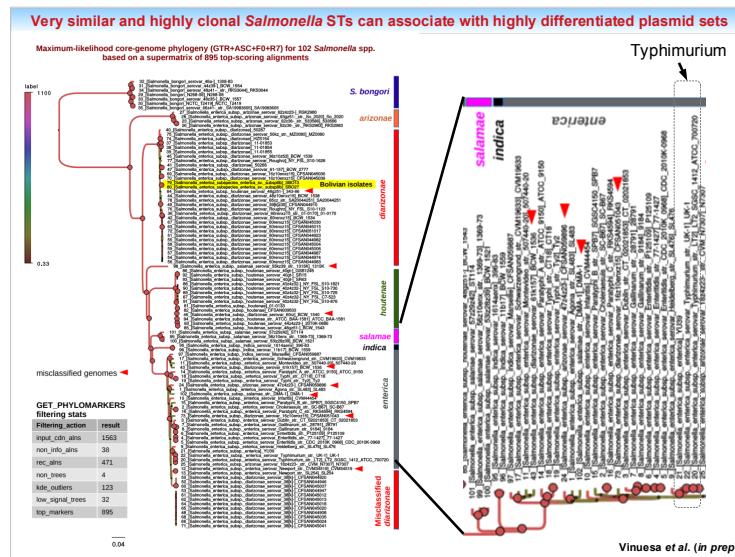
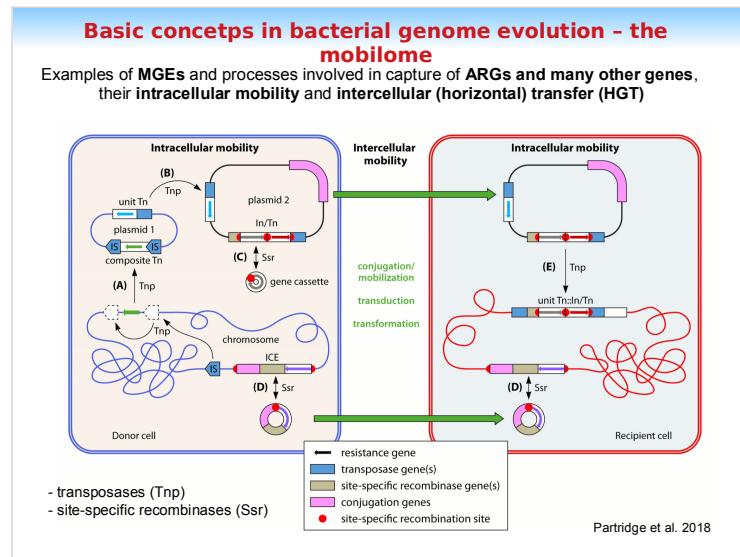
Introducción a la pan-genómica microbiana

II Introducción a la filoinformática – pan-genómica y filogenómica microbiana,
TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



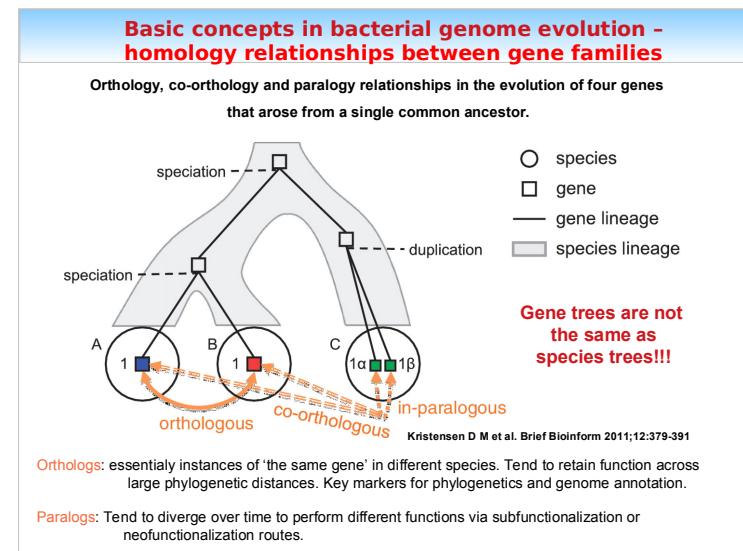
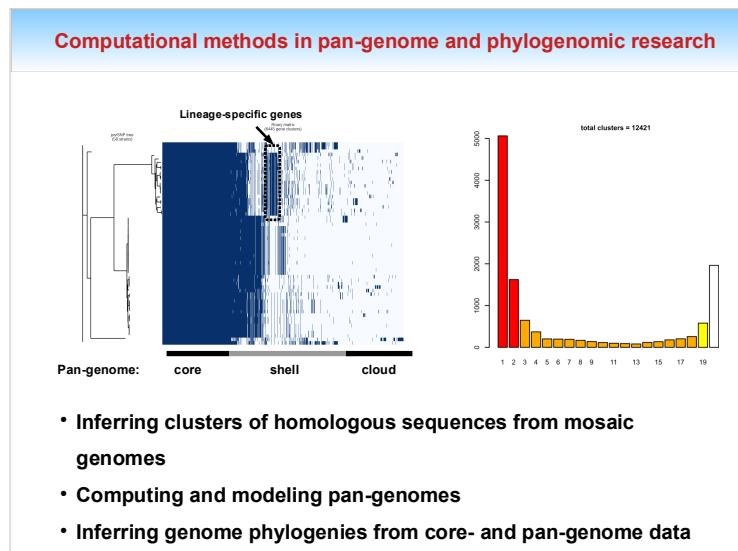
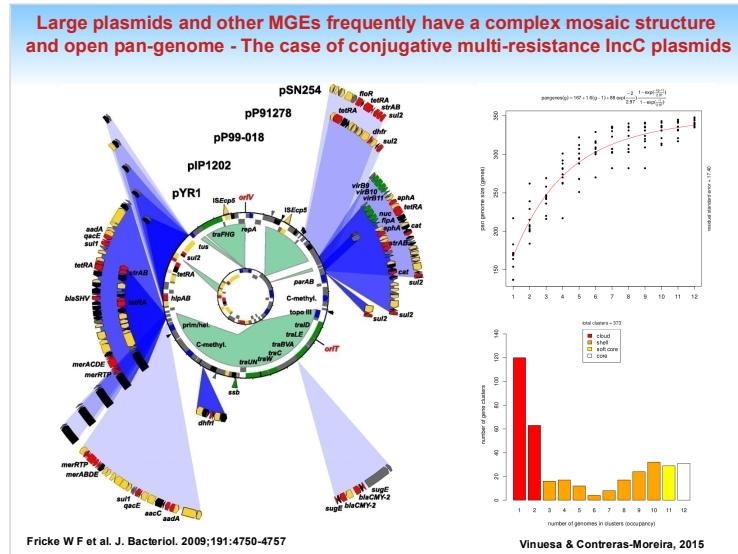
Introducción a la pan-genómica microbiana

IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana, TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



Introducción a la pan-genómica microbiana

IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana,
TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>

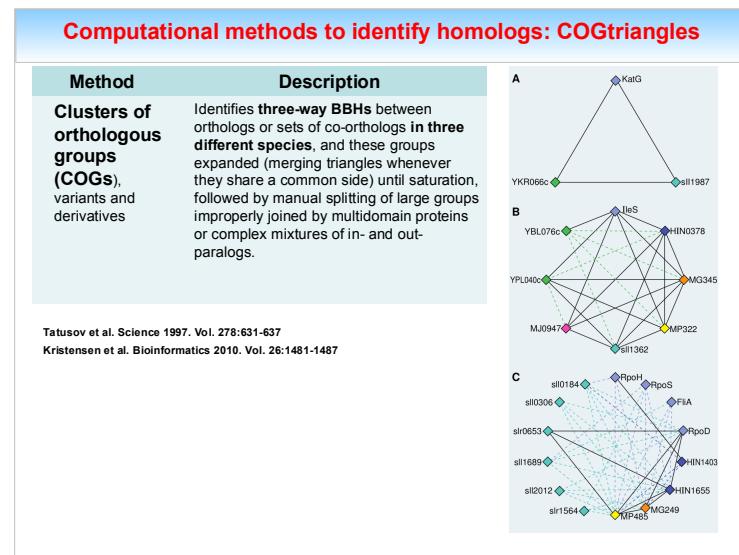
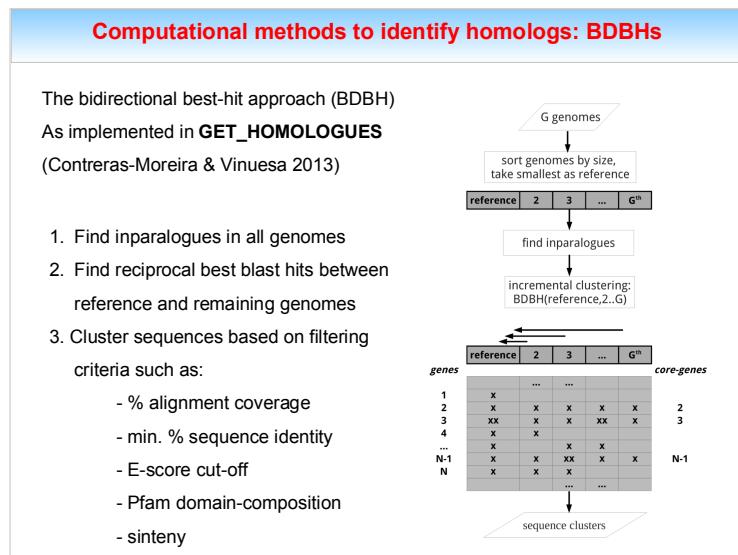
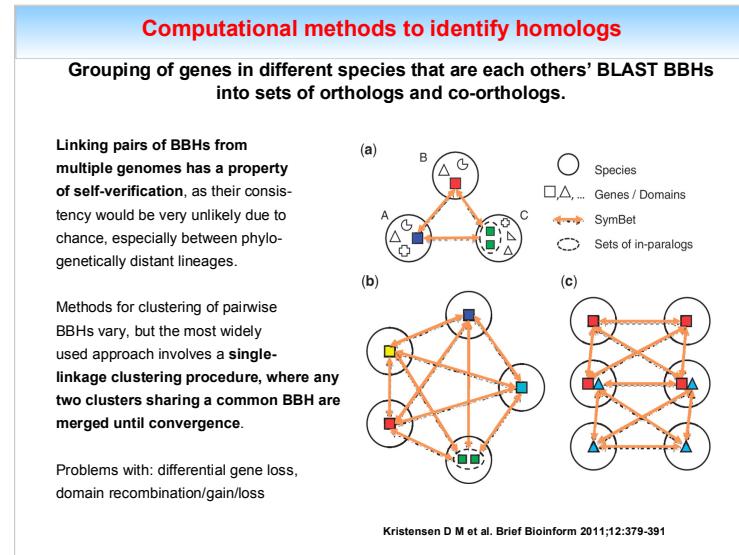
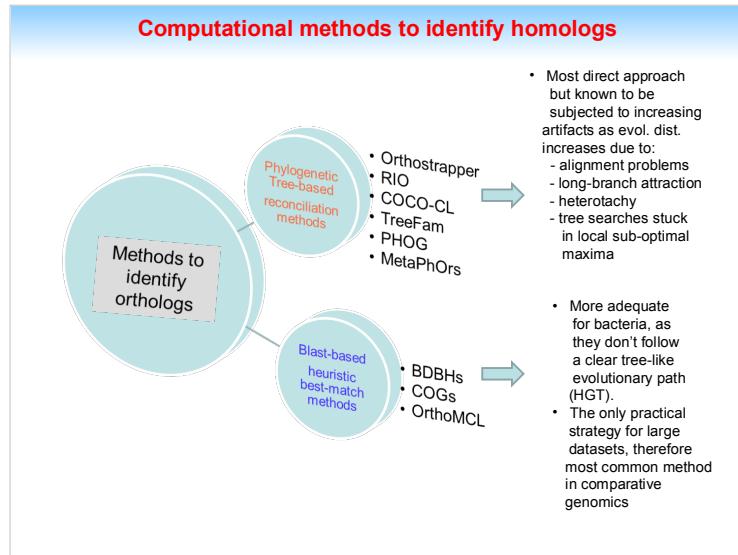


© Pablo Vinuesa 2022. @pvinmex; vinuesa{at}ccg{dot}unam{dot}mx
<https://www.ccg.unam.mx/~vinuesa/>

Licencia Creative Commons 4.0, no comercial
Attribution-NonCommercial 4.0 International (CC BY-NC 4.0)
<https://creativecommons.org/licenses/by-nc/4.0/>

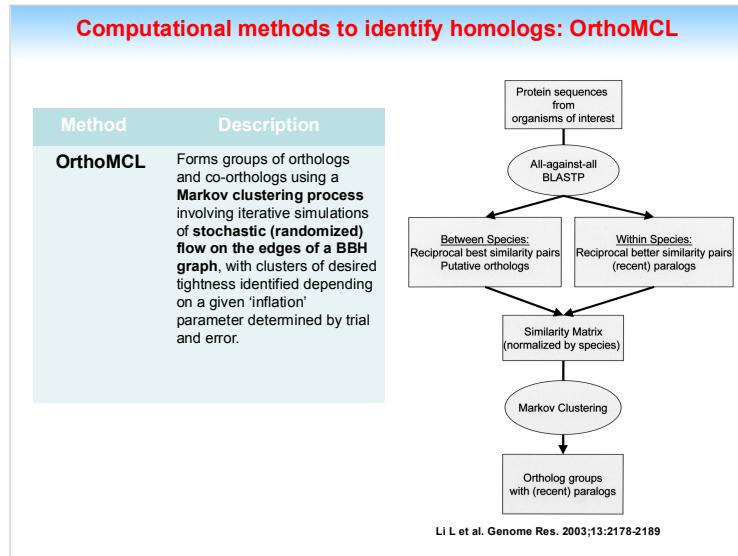
Introducción a la pan-genómica microbiana

IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana, TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



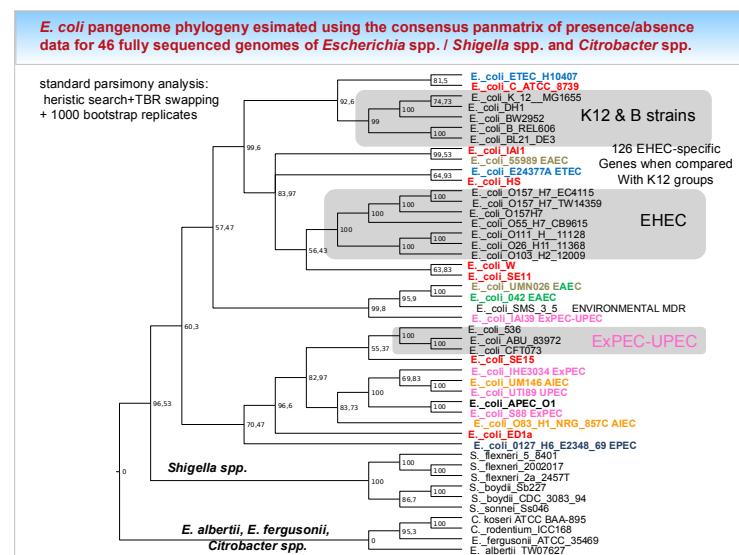
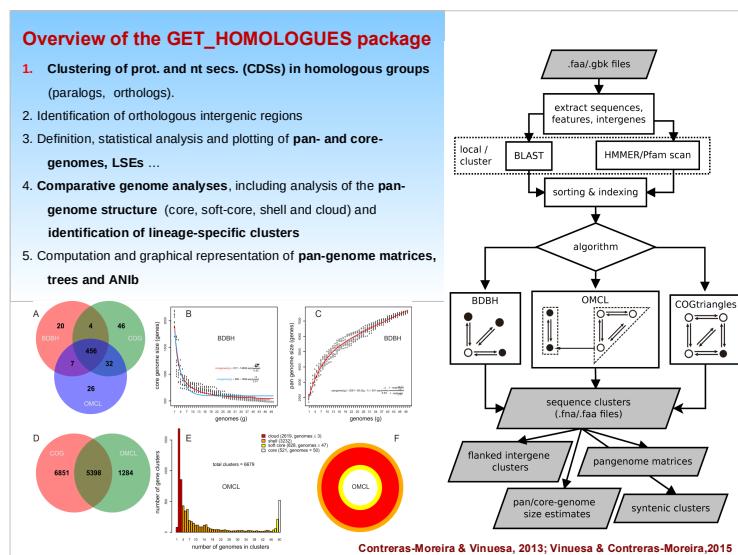
Introducción a la pan-genómica microbiana

II Introducción a la filoinformática – pan-genómica y filogenómica microbiana, TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



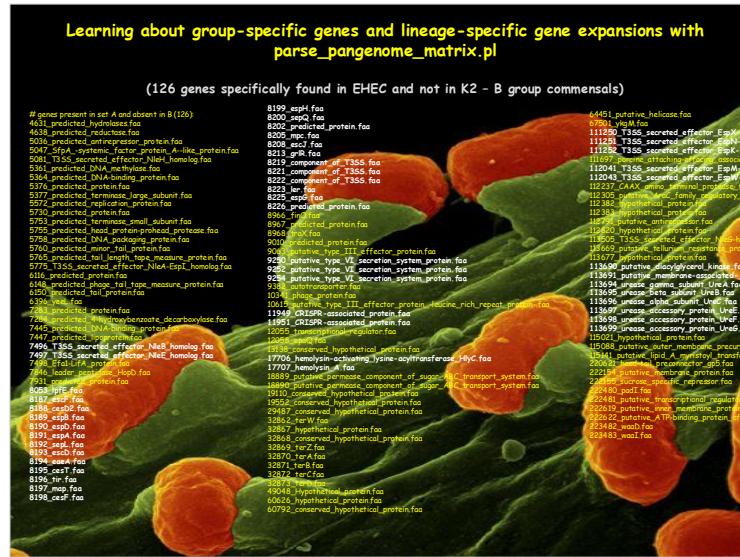
Open-source tools for phylogenomics and microbial pan-genomics GET HOMOLOGUES & GET PHYLOMARKERS

<https://github.com/vinuesa>



Introducción a la pan-genómica microbiana

IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana, TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



Introducción a la pan-genómica microbiana

IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana,
TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>

Basic Docker commands and use of containers

```
## To avoid permission errors (and the use of sudo), add your user to the docker group
# https://docs.docker.com/install/linux/linux-postinstall/
sudo groupadd docker
sudo usermod -aG docker $USER

# 1. Get general docker info and print help
$ docker info
$ docker --help
$ docker run --help

# 2. List available docker images on your system
$ docker image ls

# 3. List running containers
$ docker container ls

# 4. Stop a container
$ docker container stop CONTAINER-ID

# 5. Pull a Docker image from the registry
$ docker pull csicunam/get_homologues:latest

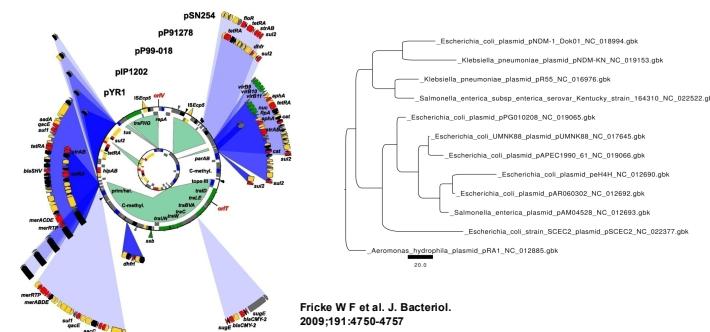
# 6. Launch an image, using a mount-bind of a user directory on the Docker container
$ docker run --rm -it -v $HOME/get_homPhy:/home/you/get_homPhy \
  csicunam/get_homologues:latest /bin/bash

# The last command uses options --rm to remove the container after exiting and sets an
-i interactive session calling a pseudo tty (-t), mounting a host directory on the
container (-v ...), accessible for wr from both, and launching a bash shell (/bin/bash)
```

Robust identification of orthologues and paralogues for microbial pan-genomics using GET_HOMOLOGUES: a case study of plncA/C plasmids

Pablo Vinuesa¹ and Bruno Condejas-Moraira^{2,3}
1 Centro de Ciencias Genómicas, Universidad Nacional Autónoma de México, Cuernavaca, Morelos, Mexico.
2 Fundación ARAID, Zaragoza, Spain.
3 Estación Experimental de Aula Dei, Consejo Superior de Investigaciones Científicas (EEAD-CSIC), Avda. Montaña, 1005, 50059 Zaragoza, Spain.

In: "Bacterial Pan-genomics". Methods in Molecular Biology. Marco Galardini, Alessio Mengoni and Marco Fondi Eds. Humana Press, Springer, 2014. *In press*



The GET_HOMOLOGUES + GET_PHYLOMARKERS tutorials:

- https://github.com/vinuesa/get_phylomarkers/

Analyses to be performed in an upcoming practical session:

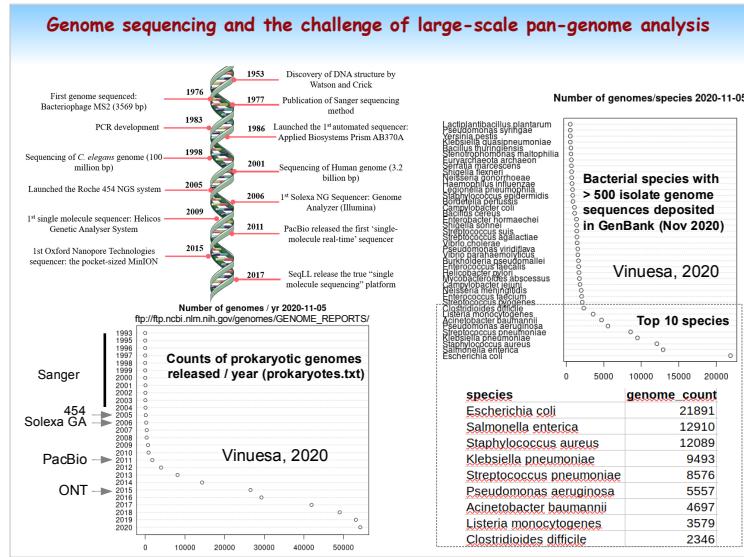
A pangenomic analysis of plncA/C plasmids using GET_HOMOLOGUES

- Defining a robust core- and pan-genome of plncA/C plasmids
- Exploring the gene space of plncA/C plasmids: core, shell, cloud
- Pan-genome trees vs. core genome trees (supermatrices)
- Identifying lineage-specific genes in NDM-1 producing plasmids

Licencia Creative Commons 4.0, no comercial
Attribution-NonCommercial 4.0 International (CC BY-NC 4.0)
<https://creativecommons.org/licenses/by-nc/4.0/>

Introducción a la pan-genómica microbiana

IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana, TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



Genome sequencing and the challenge of large-scale pan-genome analysis

Roary: rapid large-scale prokaryote pan genome analysis

Bioinformatics, 31(22), 2015, 3691–3693

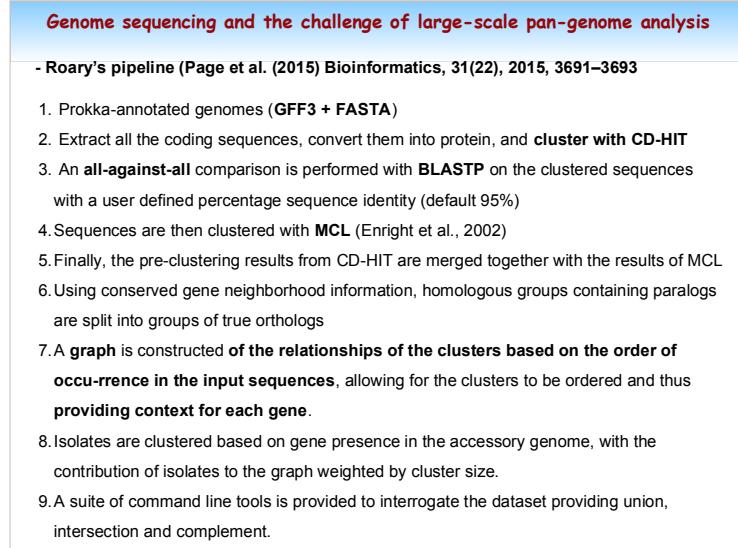
Andrew J. Page^{1,*}, Carla A. Cummins¹, Martin Hunt¹,
Vanessa K. Wong^{1,2}, Sandra Reuter², Matthew T.G. Holden³,
Maria Fookes⁴, Daniel Falush⁴, Jacqueline A. Keane³ and Julian Parkhill¹

- Standard pan-genome software approaches

- The construction of a pan genome is **NP-hard** (Nguyen et al. 2014) with additional difficulties from real data due to contamination, fragmented assemblies and poor annotation. Therefore, any approach requires **heuristics** to produce a pan genome
- all-against-all** comparison using **BLAST**, with the **running time growing ~quadratically** with the size of input data and are computationally infeasible with large datasets.
- They also have **quadratic memory requirements**, quickly exceeding the RAM available in high performance servers for large datasets.

- Roary's heuristic approach

- Roary address the computational issues by performing a **rapid clustering of highly similar sequences** with **CD-HIT** which can reduce the running time of BLAST substantially, and carefully manages RAM usage so that it increases linearly, both of which make it possible to analyze datasets with thousands of samples using commonly available computing hardware.



Genome sequencing and the challenge of large-scale pan-genome analysis

- Roary's pipeline benchmark (Page et al. (2015) Bioinformatics, 31(22), 2015, 3691–3693)

- 1 processor (AMD Opteron 6272), 60 GB of RAM.
- Constructed a simulated dataset based on *Salmonella Typhi* (S.typhi) CT18, to accurately assess the quality of the clustering.
- Created 12 genomes with 994 identical core genes and 23 accessory genes in varying combinations
- The running time and RAM of PGAP and PanOCT increases substantially with dataset size; LS-BSR, over clusters in 2% of cases.

Table 1. Accuracy of each pan genome application on a dataset of simulated data.

	Core genes	Total genes	Incorrect split	Incorrect merge
Exigent	994	1017	0	0
PGAP	991	1012	0	0
PanOCT	993	1015	1	1
LS-BSR	974	994	0	23
Roary	994	1017	0	0

Table 2. Comparison of pan genome applications using real *S. typhi* data (ERP001718)

Samples	Software	Core*	Total	RAM (mb)	Wall time (s)
8	PGAP	4345	4929	869	41 397
	PanOCT	4344	4936	663	1457
	LS-BSR	4476	4816	270	2585
	Roary	4459	4871	156	44
24	PGAP	—	—	—	—
	PanOCT	4322	4991	5313	96 093
	LS-BSR	4451	4843	554	7807
	Roary	4438	4941	444	382
1000	PGAP	—	—	—	—
	PanOCT	4272	7265	17 413	345 019
	LS-BSR	4018	9201	13 752	15 465
	Roary	4018	9201	—	—

*Core is defined as a gene being in at least 99% of samples, which allows for some assembly errors in very large datasets. Where there are no results, the application failed to complete within 5 days or used more than 60 GB of RAM. The first column is the number of unique *S. typhi* genomes in the input set with a mean of 54 contigs over all 1000 assemblies.

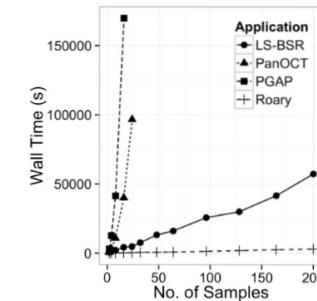
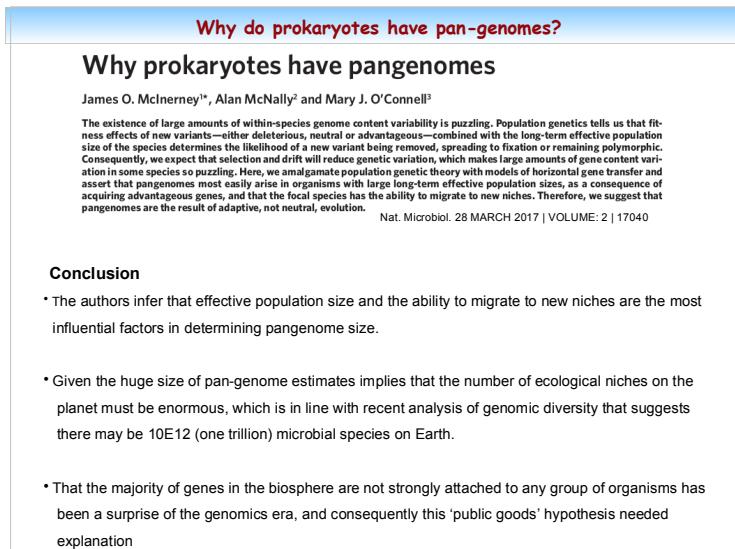
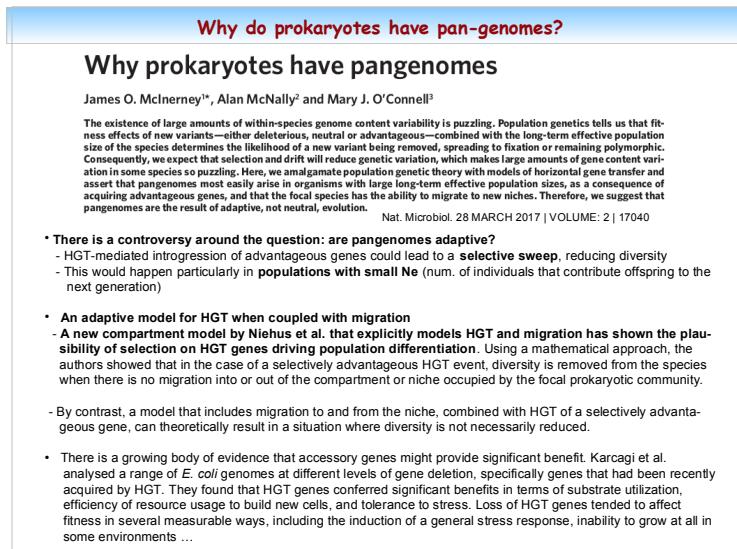
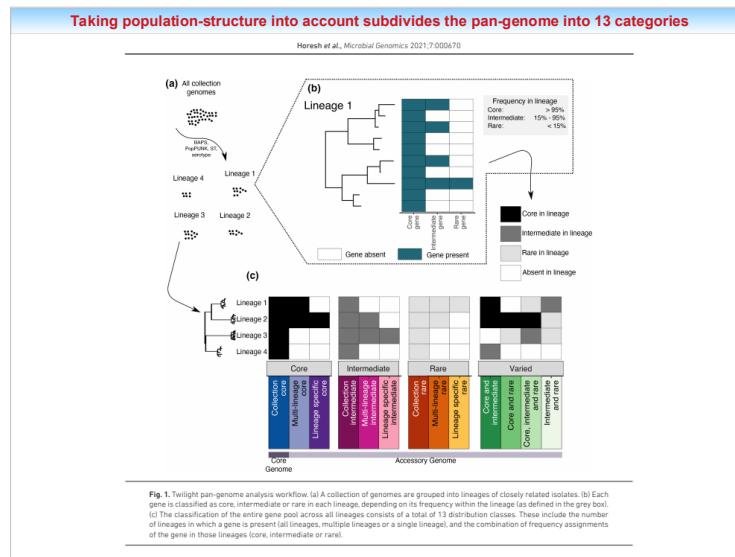
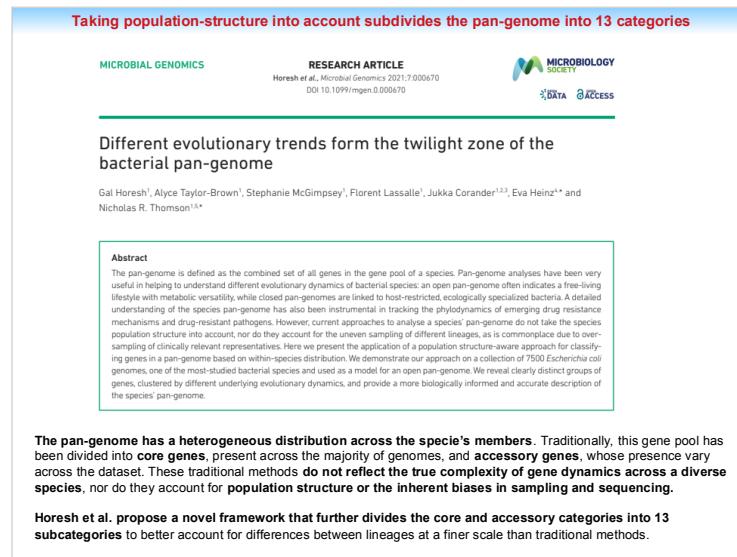


Fig. 1. Effect of dataset size on the wall time of multiple applications. On analysis that completed within 2 days and 60 GB of RAM is shown

Introducción a la pan-genómica microbiana

IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana, TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>



Introducción a la pan-genómica microbiana

IIIntroducción a la filoinformática – pan-genómica y filogenómica microbiana,
TIB2022, 1-5 Agosto, 2022 CCG-UNAM, Cuernavaca, Mor. México
<https://github.com/vinuesa/TIB-filoinfo>

Bacterial evolution, mobile genetic elements and the pan-genome

References

* Homology and genome evolution

Fitch WM. Homology a personal view on some of the problems. *Trends Genet.* 2000 May;16(5):227-31

Koonin EV. Orthologs, paralogs, and evolutionary genomics. *Annu Rev Genet.* 2005;39:309-38

Koonin EV, Wolf YI. Genomics of bacteria and archaea: the emerging dynamic view of the prokaryotic world. *Nucleic Acids Res.* 2008 Dec;36(21):6688-719.

Kristensen DM, Wolf YI, Mushegian AR, Koonin EV. Computational methods for Gene Orthology inference. *Brief Bioinform.* 2011 Sep;12(5):379-91. <=

* Bacterial mobile genetic elements

Aziz RZ, Breitbart M, Edwards RA. Transposases are the most abundant, most ubiquitous genes in nature. *Nucleic Acids Res.* 2010 Jul;38(13):4207-17

Bellanger X, Payot S, Leblond-Bouget N, Guédon G. Conjugative and mobilizable genomic islands in bacteria: evolution and diversity. *FEMS Microbiol Rev.* 2014 Jul;38(4):720-60

Cambray G, Guerout AM, Mazel D. Integrins. *Annu Rev Genet.* 2010;44:141-66

Campbell A. The future of bacteriophage biology. *Nat Rev Genet.* 2003 Jun;4(6):471-7.

Frost LS, Leplae R, Summers AO, Toussaint A. Mobile genetic elements: the agents of open source evolution. *Nat Rev Microbiol.* 2005 Sep;3(9):722-32

Galán JE, Waksman G. Protein-Injection Machines in Bacteria. *Cell.* 2018 Mar 8;172(6):1306-1318

Gillings MR. Integrins: past, present, and future. *Microbiol Mol Biol Rev.* 2014 Jun;78(2):257-77

Johnson CM, Grossman AD. Integrative and Conjugative Elements (ICEs): What They Do and How They Work. *Annu Rev Genet.* 2015;46:577-601

Juhás M, van der Meer JH, Gaillard M, Harding RM, Hood DW, Crook DW. Genomic islands: tools of bacterial horizontal gene transfer and evolution. *FEMS Microbiol Rev.* 2009

Partridge SR, Kwong SM, Firth N, Jensen SO. Mobile Genetic Elements Associated with Antimicrobial Resistance. *Clin Microbiol Rev.* 2018 Aug 1;31(4) <=

Pilla G, Tang CM. Going around in circles: virulence plasmids in enteric pathogens. *Nat Rev Microbiol.* 2018 Aug;16(8):484-495. <=

Bacterial evolution, mobile genetic elements and the pan-genome

References

* Bacterial mobile genetic elements (continuation)

Siguier P, Gourbeyre E, Varani A, Ton-Hoang B, Chandler M. Everyman's Guide to Bacterial Insertion Sequences. *Microbiol Spectr.* 2015 Apr;3(2):MDNA3-0030-2014

Welch RA, Burland V, Plunkett G 3rd, Redford P, Roesch P, Rasko D, Buckles EL, Liu SR, Boutin A, Hackett J, Stroud D, Mayhew GF, Rose DJ, Zhou S, Schwartz DC, Perna NT, Mobley HL, Donnenberg MS, Blattner FR. Extensive mosaic structure revealed by the complete genome sequence of uropathogenic Escherichia coli. *Proc Natl Acad Sci U S A.* 2002 Dec 24;99(26):17020-4

Wozniak RA, Waldor MK. Integrative and conjugative elements: mosaic mobile genetic elements enabling dynamic lateral gene flow. *Nat Rev Microbiol.* 2010 Aug 8(8):552-63.

Bacterial evolution, mobile genetic elements and the pan-genome

References

* Mosaic genome structure, population structure and the bacterial pan-genome: reviews and papers

Golicz, A. A., Bayer, P. E., Bhalia, P. L., Bailey, J. & Edwards, D. Pangenomics Comes of Age: From Bacteria to Plant and Animal Applications. *Trends in Genetics* vol. 36 132–145 (2020). <=

Horesh et al. 2021. Different evolutionary trends form the twilight zone of the bacterial pan-genome. *Microb Genom.* 2021 Sep;7(9). <=

McInerney JO, McNally MJ. Why prokaryotes have pangenomes. *Nat Microbiol.* 2017 Mar 28;2:17040 <=

Medini D, Donati C, Tettelin H, Masignani V, Rappuoli R. The microbial pan-genome. *Curr Opin Genet Dev.* 2005 Dec;15(6):589-94

Tettelin H, Masignani V, Cieslewicz MJ, Donati C, Medini D, et al. Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: implications for the microbial "pan-genome". *Proc Natl Acad Sci U S A.* 2005 Sep 27;102(39):13950-5. doi: 10.1073/pnas.0506758102. Epub 2005 Sep 19. Erratum: Proc Natl Acad Sci U S A. 2005 Nov 8;102(45):16530. PMID: 16172379; PMCID: PMC1216834.

Tettelin H, Medini D, editors. *The Pangenome: Diversity, Dynamics and Evolution of Genomes.* Cham (CH): Springer; 2020. PMID: 32633920.

Welch RA, Burland V, Plunkett G 3rd, Redford P, Roesch P, Buckles EL, Liu SR, Boutin A, Hackett J, Stroud D, Mayhew GF, Rose DJ, Zhou S, Schwartz DC, Perna NT, Mobley HL, Donnenberg MS, Blattner FR. Extensive mosaic structure revealed by the complete genome sequence of uropathogenic Escherichia coli. *Proc Natl Acad Sci U S A.* 2002 Dec 24;99(26):17020-4

* Tools and approaches for computing pan-genomes and (pan-)genome phylogenies

Contreras-Moreira B, Vinuesa P. GET_HOMOLOGUES, a versatile software package for scalable and robust microbial pangenome analysis. *Appl Environ Microbiol.* 2013 Dec;79(24):7696-701

Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, Fookes M, Flush D, Keane JA, Parkhill J. Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics.* 2015 Nov 15;31(22):3691-3.

Vinuesa P, Contreras-Moreira B. Robust identification of orthologues and paralogues for microbial pan-genomics using GET_HOMOLOGUES: a case study of pIncA/C plasmids. *Methods Mol Biol.* 2015;1231:203-232

Vinuesa P, Ochoa-Sánchez LE, Contreras-Moreira B. GET_PHYLOMARKERS, a Software Package to Select Optimal Orthologous Clusters for Phylogenomics and Inferring Pan-Genome Phylogenies, Used for a Critical Geno-Taxonomic Revision of the Genus *Stenotrophomonas*. *Front Microbiol.* 2018 May 1;9:771

Licencia Creative Commons 4.0, no comercial

Attribution-NonCommercial 4.0 International (CC BY-NC 4.0)

<https://creativecommons.org/licenses/by-nc/4.0/>